

# 1 RAID 6 Hardware Acceleration

2 MICHAEL GILROY, JAMES IRVINE, and ROBERT ATKINSON, University of Strathclyde

3 Inexpensive, reliable hard disk storage is increasingly required in both businesses and the home. As disk  
4 capacities increase and multiple drives are combined in one system the probability of multiple disk failures  
5 increases. Through the adoption of RAID 6 the capability to recover from up to two simultaneous disk  
6 failures becomes available. In this article, we present three different RAID 6 implementations each tailored  
7 to support different target applications and optimized to reduce overall hardware resource utilization. We  
8 present an optimal Reed-Solomon-based RAID 6 implementation for arrays of four disks. We also present the  
9 smallest in terms of hardware resource utilization as well having the highest throughput RAID 6 hardware  
10 solution for disk arrays of up to 15 drives. Finally, we present an implementation supporting up to 255 disks  
11 in a single array.

12 Categories and Subject Descriptors: B.3.2. [Memory Structures]: Design Styles—*Mass storage hardware*

13 General Terms: Reliability, Design, Verification

14 Additional Key Words and Phrases: RAID 6, Reed-Solomon codec

## 15 ACM Reference Format:

16 Gilroy, M., Irvine, J., and Atkinson, R. 2011. RAID 6 hardware acceleration. *ACM Trans. Embed. Comput.*  
17 *Syst.* 10, 4, Article 43 (November 2011), 17 pages.

18 DOI = 10.1145/2043662.2043667 <http://doi.acm.org/10.1145/2043662.2043667>

## 19 1. INTRODUCTION

20 The most common solution for the provision of reliable data storage is Redundant  
21 Arrays of Independent Disks (RAID) which enable multiple hard disk drives to be com-  
22 bined to offer a combination of better performance and tolerance to failure [Patterson  
23 et al. 1988]. RAID-based systems have tended to utilize high-reliability/low-density  
24 disk drives.

25 Economic storage systems must strike a delicate balance between cost, reliability,  
26 and capacity. The objective is to obtain the maximum capacity for minimum cost at a  
27 given level of reliability.

28 Traditionally, Small Computer System Interface (SCSI) disk drives were character-  
29 ized as high-reliability/high-cost devices. By contrast, ATA drives were character-  
30 ized as low-reliability/low-cost devices. The main solution for reliable RAID storage sys-  
31 tems, RAID 5, utilizes parallelism through single disk redundancy, that is, the ability  
32 to recover from a single disk failure, to increase the overall reliability of the storage  
33 system. Consequently, employing RAID storage systems based upon Advanced Tech-  
34 nology Attachment (ATA) disk technology has the potential to provide a more cost-  
35 effective alternative to comparable reliable storage systems employing SCSI or SAS

---

This work was supported by the EPSRC and A2E Ltd.

Authors' address: M. Gilroy (corresponding author), J. Irvine, and R. Atkinson, University of Stathclyde,  
Glasgow, UK: email: [mgilroy@theiet.org](mailto:mgilroy@theiet.org).

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted  
without fee provided that copies are not made or distributed for profit or commercial advantage and that  
copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights  
for components of this work owned by others than ACM must be honored. Abstracting with credit is permit-  
ted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of  
this work in other works requires prior specific permission and/or a fee. Permissions may be requested from  
the Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212)  
869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2011 ACM 1539-9087/2011/11-ART43 \$10.00

DOI 10.1145/2043662.2043667 <http://doi.acm.org/10.1145/2043662.2043667>

36 technology. The downside of RAID arrays is that the probability of more than one  
37 simultaneous disk failure increases with the number of disks in the array.

38 This represents a significant drawback in the application of low-cost/low-reliability  
39 disk drives to storage systems. By using a RAID scheme which provides mechanisms  
40 to recover data when up to two disk drives fail simultaneously a reliable, high-capacity  
41 storage system may be implemented with lower-cost disk drives.

42 The rest of this article is structured as follows: Section 2 provides an overview of  
43 RAID 5 and RAID 6. Section 3 details various RAID 6 algorithms and implementa-  
44 tions. Our initial optimal RAID 6 implementation is described in Section 4 with the  
45 reliability and performance of this design discussed in Section 5. Two improved RAID  
46 6 controller supporting arrays of up to 15 and 255 disks are described in Section 6. We  
47 then provide an implementation of the RAID 6 controller for the Altera Avalon bus in  
48 Section 7 to demonstrate the ease of use of these IP blocks and to support testing of  
49 disk arrays of up to 255 disk drives. Finally, Section 8 provides the conclusions and  
50 analysis of the article.

## 51 2. TECHNICAL BACKGROUND

52 The use of RAID was first proposed in Patterson et al. [1988] as a means of providing  
53 higher-capacity data storage and providing recovery mechanisms in the event of a  
54 single disk failure. Single large hard disk drives remain the predominant storage  
55 technology despite the known issues of complete data loss resulting from disk failure,  
56 or at the very least the need to recover data from a backup system. Backup system  
57 costs are expensive so there is a need to make the primary data storage system tolerant  
58 to disk failures.

59 The benefits of RAID include:

- 60 — protection against data loss;
- 61 — provision of real-time data recovery with uninterrupted access when a disk fails;
- 62 — multiple drives working in parallel, increasing overall system performance;
- 63 — increased system uptime and availability.

64 In addition to providing increased reliability, disk arrays improve the data throughput  
65 of storage systems [Hennessy and Patterson 2002] as a result of disk parallelism. The  
66 disadvantage associated with disk arrays is that as the number of data disks in a  
67 system increases the reliability of the system decreases, since more devices results in  
68 more chance of a failure. The Mean Time To Failure (MTTF) of a disk array is inversely  
69 proportional to the number of devices in an array, as shown in Eq. (1). We have

$$MTTF = \frac{\tau_{fail}}{D}, \quad (1)$$

70 where  $\tau_{fail}$  is the mean time to failure of a disk and  $D$  is the number of disks in the  
71 array.

72 Adding redundant storage devices to a disk array increases the array's tolerance to  
73 faults. Should a disk fail the lost information may be recovered from the redundant  
74 information stored on the remaining disks. The lost data can be successfully recovered  
75 provided further data losses, exceeding the redundancy built into the array, do not  
76 occur during the recovery period, that is, the time before the disk can be replaced and  
77 array be rebuilt.

78 The mean time to failure of modern hard disk drives is over 1 million hours [Maxtor  
79 Corporation 2005a, 2005b]. Data loss, however, is not solely dependant upon disk  
80 failure; the rate at which nonrecoverable errors occur is about 1 in  $10^{15}$ . When a 15  
81 data disk array with a single parity disk fails data lost in any one of these disk sectors  
82 due to a nonrecoverable error will result in the loss of data. Assuming a drive capacity

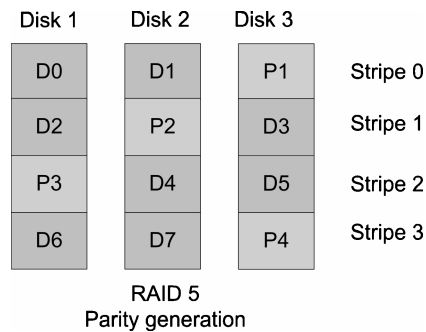


Fig. 1. RAID 5 implementation.

83 of 250GB the probability of reading all these drives successfully is 0.9758. This means  
84 that approximately 2% of all disk failures will result in data loss. Additionally, further  
85 disk failures prior to the array recovering will result in a catastrophic loss of data.

86 This article details only RAID 5 and RAID 6 implementations. A comprehensive  
87 analysis of the performance benefits of each of the most common RAID levels is avail-  
88 able in Chen et al. [1994].

## 89 2.1 RAID 5

90 RAID level 5 writes blocks of data (including parity blocks), or stripes, across disks in  
91 the array distributing the parity blocks across the array as shown in Figure 1. For a  
92 read access smaller than the block length only one disk is accessed. For larger read ac-  
93 cesses multiple disks are accessed. By distributing both the data and parity blocks the  
94 read throughput of the array is increased. This is achieved as the read/write requests  
95 are distributed evenly across each disk in the array, enabling parallel execution of disk  
96 accesses. For RAID 5 parity calculations, it is necessary to have a copy of all data in  
97 the stripe to perform the calculation; the stripe is therefore normally read into cache  
98 to maximize write throughput.

99 Figure 1 shows an example of a RAID 5 array with data stripes numbered D0-D7  
100 and the parity stripes numbered P1-P4. This pattern of data and parity blocks is  
101 distributed across the disk to improve data access speeds. Write operations which  
102 do not access all the data disks require that the parity disk be updated to reflect the  
103 changes. This is done by subtracting the modulo 2 of the difference between the old  
104 data and new data from the parity stripe. RAID 5 distributes the parity throughout  
105 the array to lower the likelihood of the parity write access becoming a bottleneck.  
106 Multiple writes are therefore possible in a RAID 5 array across different stripes.  
107 RAID 5 is the most flexible of the RAID levels, offering good support for small write  
108 operations and high read and medium write data rates.

## 109 2.2 RAID 6

110 RAID level 6 is largely undefined and was not part of the original RAID definitions  
111 proposed in Patterson et al. [1988]. It is loosely defined as offering support for double  
112 disk recovery. This is achieved through a more complex algorithm than that used  
113 in RAID 5 and requires a minimum of two checksum disks. As in RAID 5, disk  
114 striping is used to improve the write performance. An example of a RAID 6 array  
115 with two data disks and two parity disks is given in Figure 2. Here the first checksum  
116 stripe is numbered P1-P4 and the second checksum stripe numbered Q1-4. RAID 6  
117 requires both the checksum disks to be updated for every write operation decreasing

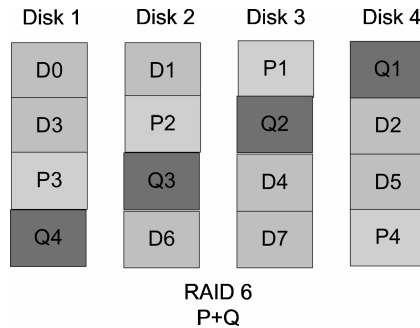


Fig. 2. RAID 6 implementation.

118 the maximum throughput achievable for write operations, but offering equivalent  
119 performance as RAID 5 for read operations.

120 RAID level 6 may be considered a solution for mission-critical applications, as well  
121 as being suited to larger disk arrays which have a higher probability of failure, espe-  
122 cially when using lower-cost, lower-reliability disk drives. As data storage capacities  
123 increase in an array the probability of read errors occurring during recovery also in-  
124 creases. With a RAID 5 implementation any read error during the rebuilding of an  
125 array will result in permanent data loss. RAID 6 is able to recover this lost data in  
126 instances where a single drive has failed. Actual data loss will not occur unless read er-  
127 rors occur across multiple disks in the stripe, or until two simultaneous drive failures  
128 occur and a read error is encountered during data recovery. During normal operation  
129 RAID 6 requires a single extra write operation when compared with RAID 5, however,  
130 read throughput matches that of a RAID 5 implementation.

131 The disadvantages of this approach are the complexity of the controller and the  
132 necessity of  $n+2$  drives to implement the design due to the two-dimensional parity  
133 scheme. RAID 6 remains nonstandardized and as such many vendor-specific RAID  
134 level 6 schemes exist. It has been suggested that standardization of RAID 6 may take  
135 place since Intel has taken an interest in RAID 6 [Karp 2007]. A number of algorithms  
136 have been proposed to support RAID 6 and these include the use of Reed-Solomon  
137 codes and orthogonal dual parity checks [Hennessy and Patterson 2002].

### 138 3. RAID 6 ALGORITHMS

139 The following section provides an overview of algorithms which have been proposed  
140 for use in RAID 6 implementations.

#### 141 3.1 EVENODD

142 EVENODD was first proposed as a suitable method for tolerating up to two disk fail-  
143 ures in a RAID system [Blaum et al. 1995]. The benefit of EVENODD is that the  
144 encoding scheme is achieved through the use of exclusive-OR operations meaning that  
145 it may be implemented through the use of a standard RAID 5 controller hardware.  
146 This scheme makes use of both horizontal and diagonal redundancy. The diagonal  
147 redundancy carries the parity information that may be either even or odd, which is  
148 where this scheme derives its name. It is optimal for two disk failures as it only re-  
149 quires two redundant disks, and small write operations only require 3 read and 3 write  
150 operations. The major disadvantage of this encoding scheme is that the write perfor-  
151 mance for certain diagonal blocks is poor, resulting in a degradation in performance of  
152 the system.

### 153 3.2 Redundancy Matrix 2 (RM2)

154 RM2 RAID 6 architectures place data and parity stripes according to a Redundancy  
155 Matrix (RM) [Park 1995]. Parity blocks are associated with groups of data disks to  
156 provide double disk redundancy. It provides the means to uniformly distribute parity  
157 data across disk arrays, across disk arrays, although by grouping disks in this manner  
158 more than two additional redundant disks may be required. The algorithm is not read-  
159 ily applicable in practice as it assumes that disks can be partitioned into any number of  
160 blocks. It is of note as it highlights the need to distribute parity data to maximize data  
161 write throughput and remove potential bottlenecks. Whilst optimizations of the basic  
162 approach have been proposed [Nam et al. 2002], RM2 remains complex to implement  
163 and is suboptimal in the number of redundant drives required.

### 164 3.3 Reed-Solomon Coding

165 Reed-Solomon (RS) codes are a class of error correcting codes which enable the recovery  
166 of any  $n$  errors from a system with  $m + n$  data inputs [Reed and Solomon 1984]. For  
167 application to RAID 6 arrays the desired property is that any 2 storage devices may fail  
168 and the data can be recovered from the remaining disks. Reed-Solomon codes provide  
169 optimal error correction capability as they require only  $n$  redundant storage locations  
170 to recover  $n$  failures.

171 Reed-Solomon-based disk arrays may operate in a manner similar to RAID 5, block-  
172 interleaved distributed-parity disk arrays. The difference is the use of Reed-Solomon  
173 coding in place of a basic parity calculation used in RAID 5 and the need for two redun-  
174 dant disks and two parity writes per stripe. Reed-Solomon encoding and decoding is  
175 performed using finite field arithmetic which has hampered RS-RAID from widespread  
176 adoption in hardware RAID controllers due to the increased complexity of the con-  
177 trollers. The benefit of a Reed-Solomon encoding scheme for RAID controllers over the  
178 other possible algorithms is that they operate over a two-dimensional array structure,  
179 unlike the three-dimensional structure of EVENODD codes, and are very similar to  
180 RAID 5 arrays. RAID 5 may be considered as a Reed-Solomon implementation for an  
181  $m + 1$  recovery system.

### 182 3.4 Galois Field Theory

183 Reed-Solomon coding is based upon Galois field theory and the a basic understanding  
184 of the theory of Galois fields is helpful in understanding the nomenclature of RS codecs.  
185 A Galois field is a field of finite order or cardinality. A Galois field of  $q$  elements is  
186 usually denoted as  $GF(q)$ . The number of elements in a finite field must be of the form  
187  $p^m$ , where  $p$  is a prime integer and  $m$  is a positive integer, and is denoted as  $GF(p^m)$   
188 or  $GF(q)$ . By giving its size, a field is described completely as for any  $q$  of the form  $p^m$ ,  
189 the field is unique up to its isomorphism: that is, only fields of the power of the prime  
190 exist.

191 The order of an element  $\alpha$  in a  $GF(q)$  is the smallest positive integer  $m$  such that  
192  $\alpha^m = 1$ . An element with order  $(q - 1)$  in  $GF(q)$  is known as a primitive element  
193 in  $GF(q)$ . There is always at least one primitive element in a  $GF(q)$ . As the  $(q - 1)$   
194 consecutive powers of  $\alpha$  must be distinct, they form the  $(q - 1)$  nonzero elements of  
195  $GF(q)$  [How 2006].

196 3.4.1 *Calculating Primitive Polynomial  $GF(4)$* . For a  $GF(4)$  the prime field is  $GF(2)$  repre-  
197 sented as binary 0 and 1. This is can be shown as  $4 = 2^2$ . All addition and multipli-  
198 cation operations performed in this field are performed modulo 2 making it ideal for  
199 hardware implementations. To find the irreducible polynomial of degree 2 in  $GF[x]$  by

200 a brute-force method, list all the quadratic rings of degree 2 with the coefficients of 0  
201 and 1.

202  $x^2$   
203  $x^2 + 1$

204 The only irreducible polynomial in this list is  $x^2 + 1$  as it does not factor to zero. This  
205 can be done again for  $GF(8)$ ,  $GF(16)$ , etc., where the field is of type  $GF(2^m)$ . Our prime  
206 will always be  $GF(2)$  and we need find the  $m$  polynomials with coefficients of 0 and 1  
207 have to be checked to find the irreducible primitives.

208 *3.4.2 Galois Field Arithmetic.* The polynomial representation of a Galois field is com-  
209 monly used for addition operations. Over a Galois field the associative and commu-  
210 tative laws apply for addition and multiplication operations. This means that both  
211 addition and subtraction are the same and may be performed by a bitwise XOR op-  
212 eration. These features allow Galois-field-based coding schemes to be readily imple-  
213 mented in hardware.

### 214 3.5 RAID 6 Implementations

215 In 2005, Intel released its IOP331 Xscale I/O processor which incorporated a hardware  
216 RAID 6 accelerator [Intel Corp. 2005]. Intel adopted a Reed-Solomon coding scheme for  
217 the RAID 6 checksum generation. The hardware acceleration units on the device assist  
218 in the calculation of parity and Galois field arithmetic. This controller is used in a  
219 number of commercial RAID controller cards to provide hardware RAID 6 solutions. In  
220 August 2006, Xilinx issued a white paper on the implementation of RAID 6 hardware  
221 acceleration for their FPGAs [Xilinx Inc. 2006].

222 Reed-Solomon codes appear to be the solution of preference for commercial hard-  
223 ware RAID 6 controllers at present. Software RAID 6 based upon a Reed-Solomon  
224 encoding scheme is also available as part of the standard Linux kernel and has been  
225 included since version 2.6 was released. This allowed for the hardware implemen-  
226 tations that we describe shortly to be directly compared with a comparable software  
227 implementation.

## 228 4. THE DESIGN OF A RAID 6 CONTROLLER

229 In this section we present the first of our hardware RAID 6 implementations. This  
230 implementation is optimal for arrays of four disks and was designed to require minimal  
231 hardware resources. Whilst RAID 6 is not commonly used in such small disk arrays  
232 there are applications, such as CCTV recording and other security applications, where  
233 the ability to remove two drives and generate a direct replacement of all the data is of  
234 benefit.

235 All FPGA development presented this section was performed using Verilog HDL  
236 and simulated in Mentor Graphics ModelSim PE. Where possible system-level test-  
237 benches were written, again in Verilog HDL, to simulate the hardware design and  
238 also verify through built-in checks that the hardware operated as expected. The test-  
239 bench verified the accelerators, ability to encode and decode all combinations of in-  
240 put data for all possible array combinations. Synthesis was undertaken using the  
241 Altera Quartus II design tools. Our PCI-based hardware was implemented on an Al-  
242 tera Stratix EP1S20F780 PCI-X/PCI Development Kit 32/64-bit from PLDA [PLDA  
243 2002]. All other hardware implementations were implemented on an Altera Cyclone  
244 II development board.

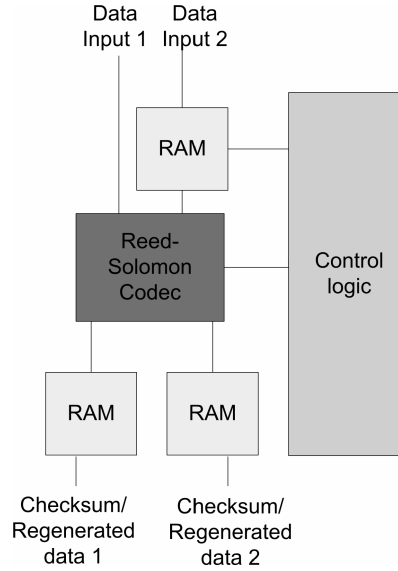


Fig. 3. Hardware implementation of the Reed-Solomon encoder using  $GF(4)$ .

#### 245 4.1 Reed-Solomon Codec

246 The codec presented in this section operates in a  $GF(4)$  and is the basis of an optimal  
 247 RS-RAID 6 accelerator for arrays of four storage devices. The codec is optimal for this  
 248 purpose as it can support only two data disks and two checksum disks.

249 The hardware codec is designed to operate under the following assumptions.

- 250 (1) The data disks and checksum disks are located at fixed locations, that is, data disks  
 251 at location 1 and 2 and checksum at location 3 and 4.  
 252 (2) Both data inputs are passed into the codec simultaneously.  
 253 (3) 2 bits per input are encoded or decoded simultaneously.

254 The checksum is generated according to a generator function. Encoding of the first  
 255 checksum is performed by a simple bitwise XOR operation of the two data inputs, both  
 256 two bits wide. Generation of the second checksum required that each input be multi-  
 257 plied by a fixed value over the Galois field. For this codec the first input is multiplied  
 258 by one and the second by two. These multiplicative operations are performed by a  
 259 lookup table. The hardware implementation of the Reed-Solomon encoder is shown in  
 260 Figure 3.

261 Figure 4 shows the path through the decoding hardware when both of the erasures  
 262 occur in the data disks.

263 If one or both of the checksum disks are erased then the lost data is regenerated  
 264 by passing the data through the Reed-Solomon encoder and selecting the appropriate  
 265 output. If one or both of the erased disks contained data then the missing data is  
 266 regenerated using the appropriate decoding functions as described in Eqs. (2) and (3).

$$\begin{aligned}
 (d_i \times i) \oplus (d_j \times j) &= c_1 \oplus s_2 \times i \\
 \Rightarrow (d_i \times i) \oplus (d_j \times j) &= (c_2 \oplus s_2) \oplus [(c_1 \oplus s_1) \times i]
 \end{aligned} \tag{2}$$

267

$$(d_i \times i) \oplus (d_j \times j) = c_1 \oplus s_2 \Rightarrow \frac{(c_2 \oplus s_2) \oplus [(c_1 \oplus s_1) \times i]}{i \oplus j} \tag{3}$$

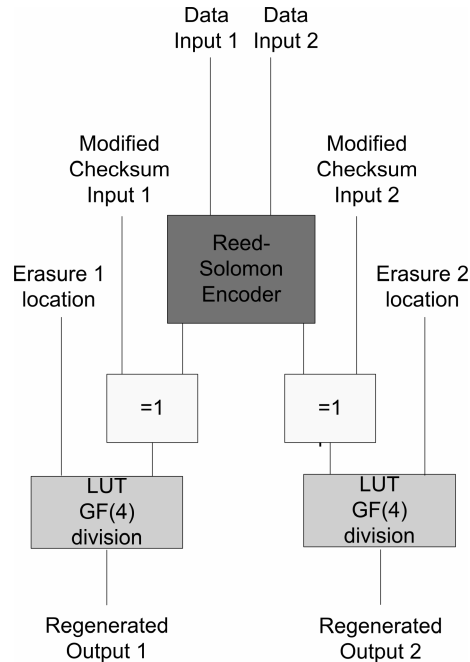


Fig. 4. Hardware implementation for the Reed-Solomon decoder using  $GF(4)$  for double data disk recovery.

268 The RS codec operation was verified by encoding all combinations of two-bit inputs  
 269 and storing the results. All single and double disk failures are tested using the stored  
 270 data and verifying the output with the original data. This exhaustive test of the codec  
 271 demonstrates that the encoder and decoder operate correctly for all possible data in-  
 272 puts, and recover from all permutations of single and double erasures.

#### 273 4.2 RS RAID Hardware Accelerator

274 To demonstrate this RAID 6 implementation in a practical system and verify its per-  
 275 formance a PCI interface was added to our design. This necessitated further control  
 276 logic and packaging of our original implementation. Instantiating multiple instances  
 277 of the codec in parallel allowed us to match the codec data input width to that of the  
 278 PCI bus width. To control data flow and ensure data is processed correctly a memory  
 279 interface is added to allow data to be buffered before and after encoding/decoding takes  
 280 place. Additional control logic to control the flow of data into and out of the codecs and  
 281 configure the encoding and decoding operations is also provided. A block diagram of  
 282 the RS RAID accelerator is shown in Figure 5.

283 The memory buffers are implemented to enable the first block of input data from  
 284 the first disk to be stored locally within the RAID accelerator block. The second disk's  
 285 data is input to the codec whilst simultaneously reading the buffered data.

286 The hardware was tested by simulation for all combinations of double and single  
 287 disk failures using the previously generated Verilog HDL testbenches. The RS codec  
 288 was tested by encoding all combinations of two-bit inputs and storing the results. All  
 289 single and double disk failures were then tested one at a time using the stored data and  
 290 verifying the output with the original data. These showed the hardware accelerator  
 291 was able to correctly encode and recover all combinations of double and single disc



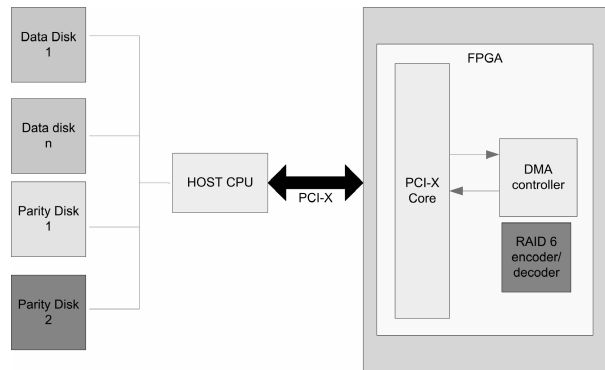


Fig. 5. RS RAID accelerator block diagram.

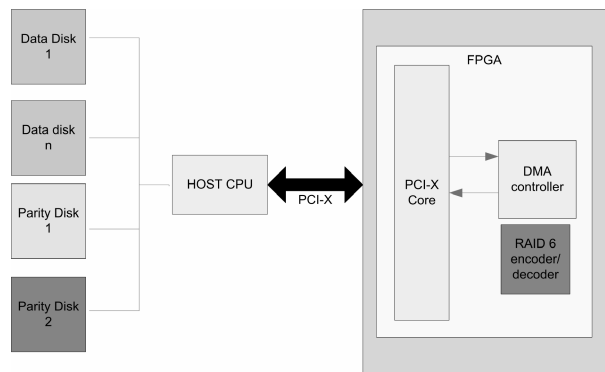


Fig. 6. PCI-based RAID 6 accelerator block diagram.

292 erasure. The results of timing analysis after synthesis for a Altera Cyclone II showed  
 293 the accelerator capable of operating at a clock rate in excess of 200 MHz.

### 294 4.3 PCI RAID 6 Accelerator System Architecture

295 Having developed the RS RAID accelerator it was necessary to quantify the perfor-  
 296 mance benefit achieved through the use of hardware acceleration. The accelerator  
 297 block was implemented on an Altera Stratix FPGA-based PCI-X development board  
 298 with a PCI 64-bit 66 MHz bus, DMA controller, and memory blocks. To enable rapid  
 299 software driver development the software RAID 6 implementation of the Linux 2.6 ker-  
 300 nel was modified to support the hardware accelerator. All comparisons in performance  
 301 between the software RAID 6 and our hardware implementation were made using the  
 302 unmodified kernel.

303 The test platform consisted of: a server class computer running Linux kernel 2.6.9.2  
 304 on an AMD Opteron processor with 1GB of RAM; 1 parallel ATA drive for the operating  
 305 system; and four 200GB Maxtor DiamondMax Serial ATA drives for testing RAID 6  
 306 arrays attached via a four port SATA controller, Si3114, from Silicon Image [Maxtor  
 307 Corporation 2005b; Silicon Image 2007].

308 To determine the read/write throughput during normal disk operations Bonnie++  
 309 [Bray 2006] was used to provide benchmark results for the test platform. Bonnie++  
 310 is a benchmarking tool used to determine data throughput performance for different  
 311 types of read and write accesses on hard disk drives. To ensure that all the required

Table I. Benchmark Results for Test Platform Using Software RAID and Hardware RAID 6

Array Configuration	Sequential read throughput	Sequential write throughput
Single Disk	50.3 MB/sec	49.7 MB/sec
RAID 0	209 MB/sec	197 MB/sec
RAID 5	102.8 MB/sec	90.4 MB/sec
software RAID 6	68.4 MB/sec	59 MB/sec
Hardware RAID 6	59.5 MB/sec	59.3 MB/sec

312 data could not be cached, which may result in unrealistic performance metrics, a 4GB  
 313 file size was used for read and write operations during benchmarking. A summary  
 314 of the results of this benchmarking is given in Table I. All benchmark figures were  
 315 generated using the unmodified Linux kernel and software RAID implementations.

316 Our test platform had a number of limitations for testing our RAID 6 accelerator.  
 317 As we did not have direct attachment to the hard disk drives, data must be read from  
 318 the drives by the software driver then passed to the RAID 6 hardware accelerator. The  
 319 encoded/decoded data was then passed back to the CPU to then be written back to  
 320 the hard drives. As a direct result of this limitation the maximum data throughput  
 321 when using the hardware accelerator will be halved. Also, as the software driver for  
 322 controlling access to the disk drives did not reside on the hardware accelerator we  
 323 expect a further performance degradation compared to a complete stand-alone RAID 6  
 324 hardware controller.

## 325 5. RELIABILITY TESTING

### 326 5.1 Array Generation and Recovery

327 To verify that the hardware accelerator operated correctly within the test system  
 328 a RAID 6 array was generated and verified using the Linux *mdadm* tools and the  
 329 modified RAID 6 device driver. Using the hardware accelerator it was possible to  
 330 successfully:

- 331 — generate a RAID 6 array;
- 332 — build a Linux file system;
- 333 — mount the RAID 6 array;
- 334 — read and write files to the array.

335 Verification of data write integrity has been conducted by comparing a known written  
 336 file with its original and checking for any differences. In this and all other tests each  
 337 of the four disk drives utilized 10GB partitions, offering 20GB of usable data storage  
 338 space. The 10GB partitions were specifically selected to ensure that the test results  
 339 were not skewed in favor of the software-only solution due to the 1GB of RAM available  
 340 on the server and to enable the time taken to build and rebuild arrays to be kept to  
 341 less than one hour.

342 With the RAID 6 array functioning correctly under normal operating conditions,  
 343 an exhaustive test was carried out of all possible combinations of disk failure and  
 344 recovery to verify the hardware simulation results. Each disk was removed and re-  
 345 built separately; files are written to and read from the array whilst a single disk is  
 346 removed (operating in the degraded mode) and during the rebuild stage. Once the  
 347 array is rebuilt, the files on the array are compared with the original files to verify  
 348 that the data has been successfully recovered. These steps are repeated with each  
 349 permutation of double disk erasure. Again the array successfully regenerated the  
 350 missing data.

Table II. Data Throughput for Software and 64-bit 66 MHz PCI Bus-Based Hardware Accelerator

	Software RAID 6	Hardware RAID 6
Array generation speed	17 MB/sec	23 MB/sec
Data recovery 1 erasure	17 MB/sec	23 MB/sec
Data recovery 2 erasure	15 MB/sec	25 MB/sec

Table III. Bonnie++ Results for Sequential Read and Write Operations Using the Hardware Raid 6 Accelerator and the Software RAID Implementation

	<i>Sequential block read operations</i>		<i>Sequential block write operations</i>	
	Throughput KB/sec	% CPU	Throughput KB/sec	% CPU
HW normal operation	59507	25	59347	7
SW normal operation	68427	14	59098	7
HW single erasure	29334	7	74910	10
SW single erasure	30734	7	92891	12
HW double erasure	33527	8	76531	12
SW double erasure	32884	7	67390	8

351 Direct comparison of the software RAID 6 solution with the hardware accelerated  
 352 solution is shown in Table II. The performance benefit of the hardware-accelerator is  
 353 not as significant as would be expected due to system architecture of having to access  
 354 the hard disk drives via the PCI bus as discussed in Section 4.3.

### 355 5.2 Read/Write Performance

356 The read/write performance of both the hardware and software RAID 6 were bench-  
 357 marked using Bonnie++. There was a clear improvement shown in using the hard-  
 358 ware accelerator over the software-only solution when a double disk failure had  
 359 occurred. A short extract of the results for both software and hardware RAID 6 is  
 360 presented in Table III. It should be noted that the software RAID 6 was operat-  
 361 ing under ideal operating conditions. The hardware-accelerated system was using  
 362 an inefficient architecture as the hard disk drives were not directly connected to the  
 363 FPGA.

364 These additional bus transactions result in the hardware accelerator operating at  
 365 least half as efficiently as would be expected in a full RAID 6 controller implementa-  
 366 tion. When migrated to a more appropriate RAID controller design it would be ex-  
 367 pected that the accelerator will provide at least twice the data throughput as that  
 368 shown here.

369 The software driver performance was found inconsistent with the data throughput,  
 370 dropping dramatically as disks failed, whereas the hardware performed consistently  
 371 during all modes of operation. It was also noted the software was prone to periodic  
 372 drops in performance for short periods which resulted in low data throughput rates.  
 373 These drops in data throughput corresponded to increases in the CPU utilization of  
 374 the test systems, background processes and would appear to demonstrate that the  
 375 software performance was highly dependant upon the loading of the server. These  
 376 drops in performance could be by as much as 90% of the average data throughput  
 377 rates. This was not observed when using the hardware accelerator.

378 The hardware accelerator was able to produce a data throughput rate which  
 379 closely matched those of the software solution under normal operation, offering a

380 slight gain during single disk failures and a greater benefit during double disk  
381 erasures.

### 382 5.3 CPU Utilization

383 The hardware accelerator appeared to be occupying large amount of CPU time  
384 when performing the Bonnie++ benchmarking tests. This was the opposite of what  
385 was expected. Upon closer examination of the processes running on the CPU dur-  
386 ing the testing it was found that this CPU utilization was actually being used by  
387 Bonnie++ and not the hardware device driver. During the software RAID 6 test-  
388 ing the majority of the CPU was utilized performing RAID 6 parity calculations  
389 and accesses. During hardware testing the majority of the CPU time was spent  
390 by Bonnie++ continually checking if the data had been processed or not. The ac-  
391 tual CPU utilization of the hardware driver was found to be around 3–4% during  
392 all operational modes. The software driver CPU utilization was found to be around  
393 33% higher during normal operation increasing to over 50% higher during double  
394 disk failures. These results may indicate that Bonnie++ was not the ideal tool to  
395 use to perform our benchmarking, however, it was the best solution available at  
396 the time.

## 397 6. HARDWARE RAID 6 FOR UP TO 15 DISK ARRAYS AND BEYOND

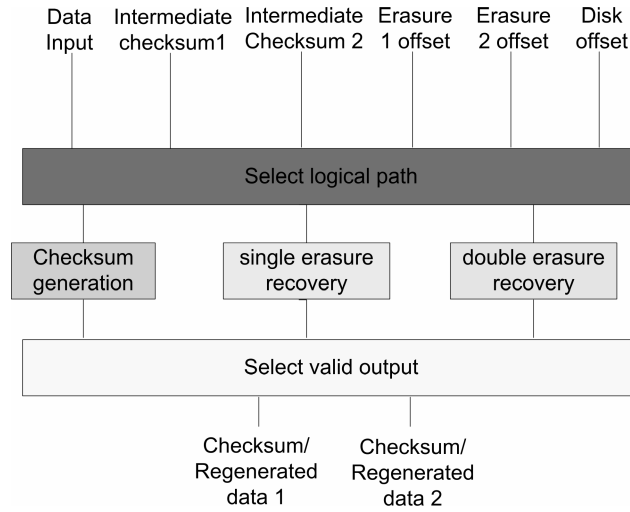
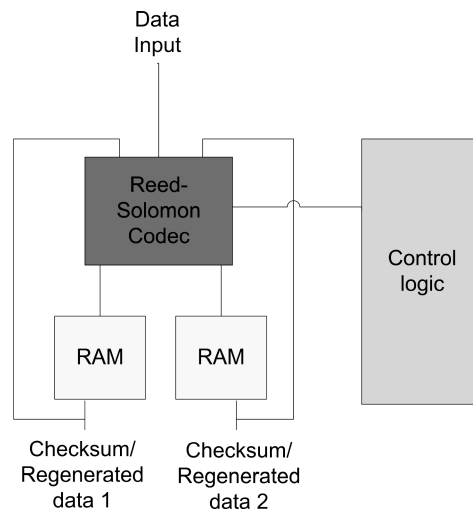
398 Having successfully demonstrated the hardware RAID 6 accelerator using a Reed-  
399 Solomon codec using  $GF(4)$  arithmetic, it was necessary to produce an improved ver-  
400 sion which offered support for a larger number of disk drives. Two implementations  
401 are presented in this section. The first provides support for up to 15 disks in a single  
402 array and is the RS coding performed over a Galois Field of 16,  $GF(16)$ , using 4-bit  
403 wide inputs as opposed to the 2-bit wide inputs of the original codec. The target ap-  
404 plication in this instance was to support a stand-alone storage system providing dual  
405 disk redundancy. Larger disk arrays of up to 255 drives are supported when perform-  
406 ing calculations over a Galois field of 256,  $GF(256)$ . In this instance a input bus width  
407 of 8 bits is used.

408 Where the original RS codec provided separate encoder and decoder blocks oper-  
409 ating within strict limits for input orders, the new codecs were designed to be more  
410 flexible and accommodating. The codecs required the following additional information:

- 411 — number of disks in the array;
- 412 — offset of the first checksum;
- 413 — offset of the second checksum;
- 414 — current disk offset;
- 415 — location of disk erasures;
- 416 — modified checksum input data.

417 By providing this information the codecs are able to perform the checksum generation  
418 or data recovery automatically as the data path through the codec is altered based  
419 upon these inputs.

420 The improved codecs, instead of requiring all the input data to be buffered prior to  
421 performing the coding, encodes the input immediately. This requires that the results of  
422 previous disk inputs, the modified checksum, be input along with the new data input.  
423 A simplified block diagram showing this improved codec design for the is shown in  
424 Figure 7.

Fig. 7.  $GF(16)$  or  $GF(256)$  codec.Fig. 8.  $GF(16)$  RS RAID 6 hardware accelerator.

#### 425 6.1 Modified RAID 6 Accelerator

426 To enable generation of the appropriate checksum for larger arrays when data is pro-  
 427 cessed on a disk-by-disk basis the intermediate checksum which is generated as each  
 428 individual disk is processed must be stored for use in processing the next disk in the  
 429 stripe. To ensure data throughput performance this intermediate checksum was stored  
 430 locally and this memory block was configured to automatically output the previously  
 431 encoded data into the codec as new data arrived, as illustrated in Figure 8. This design  
 432 change freed up memory resources consumed by the previous RAID 6 accelerator used  
 433 to buffer the input prior to performing coding operations. Again multiple instances of  
 434 the codec were instantiated in parallel to support 32-bit and 64-bit coding operations.

Table IV. Resource Utilization of the Hardware RAID 6 Accelerator Blocks

	$GF(4)$ codec	$GF(16)$ codec	$GF(256)$ codec
LUT Utilisation per RS codec	63	304	1243
LUT Utilisation for RS RAID 6 accelerator	841	4854	7105
Maximum clock rates on a Stratix FPGA	204.7 MHz	198.7 MHz	197.6 MHz
Coding latency	3 clock cycles	7 clock cycles	8 clock cycles

## 435 6.2 Comparison of RAID 6 Accelerators

436 The  $GF(4)$ -based Reed-Solomon codec has been optimized to take advantage of a fixed  
 437 structure for data to be input and output. This enables minimization of the logic re-  
 438 sources required for the  $GF(4)$  codec. To reduce the overheads associated with config-  
 439 uring the  $GF(16)$  and  $GF(256)$  codecs and to make it more flexible it has been designed  
 440 to be agnostic to: the input sequence, the number of disks in the array (provided there  
 441 are no more than 15), and offer greater flexibility than available from the original  
 442 codec. The overall resource utilization and maximum clock rates for both codecs are  
 443 detailed in Table IV.

444 Despite the  $GF(16)$  codec being designed for flexibility rather than optimized for  
 445 hardware resource consumption like the original codec, the resource utilization of the  
 446  $GF(16)$  codec was 20% of that demonstrated by the only other known published imple-  
 447 mentation [Song et al. 2005]. The codec may also be clocked at a higher rate with a  
 448 lower clock latency than this area-efficient architecture.

449 The modified RS RAID 6 hardware accelerators were added seamlessly into the  
 450 existing PCI-based test platform. The software driver was modified to provide the  
 451 additional configuration data required by these codecs.

452 An additional Silicon Image 3114 SATA controller card was added to provide access  
 453 to a fifth SATA disk drive. Verification of a five-disk RAID 6 array was also performed  
 454 at this point. The differences in performance for the hardware-accelerated design and  
 455 the unmodified software RAID 6 solution were equivalent to those found using an  
 456 array of four disks. This shows that the performance of the hardware RAID 6 codec  
 457 continues to perform as expected.

458 Larger disk arrays were not benchmarked due to a lack of suitable SATA disk drives  
 459 and a lack of space in the test system to add additional drives. However, by partitioning  
 460 the five disk drives which were available, larger arrays were shown to work with the  
 461 modified RAID 6 hardware accelerator.

## 462 7. RAID 6 IP BLOCK FOR AVALON BUS

463 To enable rapid verification of the  $GF(16)$  and  $GF(256)$  codecs in hardware an Avalon  
 464 bus wrapper was added to the design. The benefits of implementing the RAID 6 accel-  
 465 erator on the Avalon architecture was that it enabled testing and performance metrics  
 466 to be gathered on the RAID 6 accelerator for all combinations of disk array. For a  
 467 more generic implementation the PCI control was stripped from the DMA controller  
 468 and an Avalon bus wrapper was added to the IP block. The Avalon bus based RAID 6  
 469 accelerator is shown in Figure 9.

### 470 7.1 Avalon Test System

471 A testbench was used to verify basic Avalon bus transactions. However, as the  
 472 SOPC builder enabled rapid system design and development it was decided that

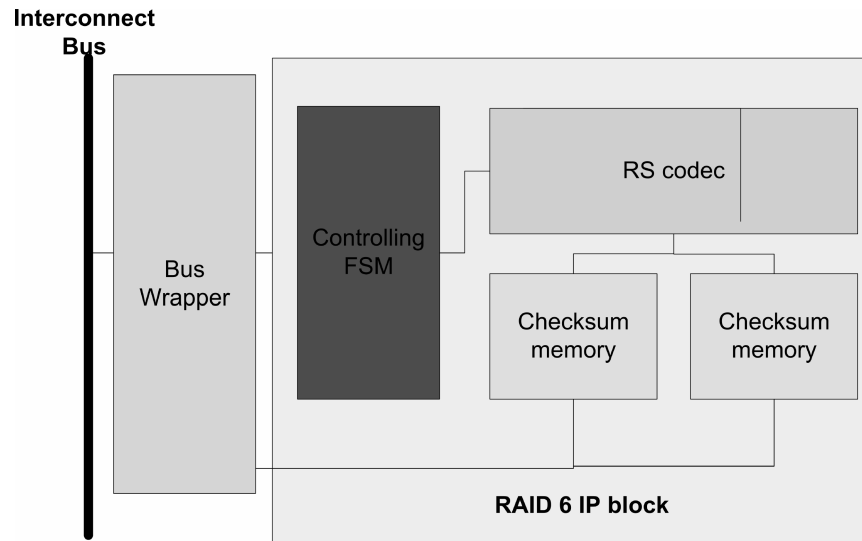


Fig. 9. Avalon bus-based RAID 6 IP block.

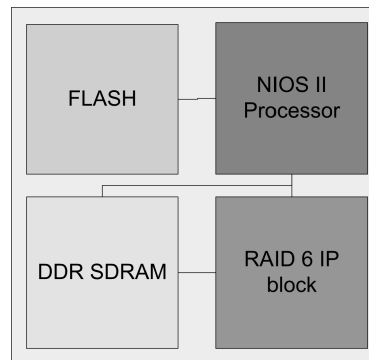


Fig. 10. Avalon bus-based test system on the cyclone II development board.

473 the generation of a simple Avalon bus-based system under the control of a NIOS  
 474 II processor would produce more rapid and valid test results. Using the Cyclone II  
 475 development kit an example system was set up with access to the following.

- 476 — NIOS II processor: Provides overall system control and a simple and rapid means  
 477 to generate test data and verify the returned data.
- 478 — Serial EEPROM: Used to store the software programs for the NIOS.
- 479 — DDR SDRAM interface: Connected to 64MB of DDR RAM to be used as program  
 480 memory.
- 481 — RAID 6 IP block.

482 This hardware configuration is shown in Figure 10. These four blocks were con-  
 483 nected via the Avalon bus using the SOPC builder and a simple test system was  
 484 synthesized.

485 Software was developed for the NIOS II processor using the Alteras NIOS II  
 486 development kit with the source code written in C. A program was developed which  
 487 segmented regions of the DDR RAM to act as four disks in a RAID 6 array. Data was

488 written to two of these disks and the RAID 6 hardware accelerator was configured to  
489 generate the checksums for the remainder. Once the RAID 6 array was successfully  
490 generated, the data disks were marked as having suffered an erasure and the data  
491 regenerated from the checksum. This regenerated data was checked with the original  
492 data to verify that the data had been recovered correctly.

493 This testing was repeated for all combinations of single and double disk erasure and  
494 for larger disk arrays. In every instance the hardware RAID 6 accelerator successfully  
495 regenerated the lost data and wrote it back to the correct locations.

## 496 8. CONCLUSION

497 We have developed and demonstrated three Reed-Solomon-based RAID 6 controllers,  
498 each designed to address separate target applications. Our RAID 6 accelerators are  
499 readily interchangeable to enable the optimal implementation to be selected for a given  
500 application. The first provided an optimal implementation for arrays of four disks. The  
501 second RAID 6 accelerator provides support for up to 15 disks in an array. The results  
502 of testing this design has shown that it requires only 20% of the physical resources  
503 of the previous best implementation as well as operating at higher clock rate. This  
504 enables our hardware accelerator to provide a higher data throughput than previously  
505 shown. Our final implementation provided support for the more common  $GF(256)$ -  
506 based RS RAID 6 controller. The scalability of the  $GF(256)$  codec has made it the  
507 implementation of choice for commercial hardware RAID 6 implementations. From our  
508 investigation we have demonstrated that significant reductions in hardware resources  
509 may be made through the adoption of the algorithm most appropriate to the target  
510 application. At the same time we have sustained the same data throughput for each  
511 separate codec.

## 512 ACKNOWLEDGMENTS

513 The authors would like to thank the EPSRC for their financial support. They would also like to thank A2E  
514 Ltd. for their financial and technical support throughout the project.

## 515 REFERENCES

- 516 BLAUM, M., BRADY, J., BRUCK, J., AND MENON, J. 1995. Evenodd: An efficient scheme for tolerating double  
517 disk failures in raid architectures. *IEEE Trans. Comput.* 44, 2, 192–202.
- 518 BRAY, T. 2006. *BONNIE++*. <http://code.google.com/p/bonnie-64/>.
- 519 CHEN, P. M., LEE, E. K., GIBSON, G. A., KATZ, R. H., AND PATTERSON, D. A. 1994. RAID: High-  
520 Performance, reliable secondary storage. *ACM Comput. Surv.* 26, 2, 145–185.
- 521 HENNESSY, J. L. AND PATTERSON, D. A. 2002. *Computer Architecture: A Quantitative Approach*. Morgan  
522 Kaufmann Publishers, San Francisco, CA.
- 523 HOW. 2006. *Fields and Galois Theory*. Springer.
- 524 INTEL CORP. 2005. *Intel 80331 IO Processor Datasheet*. Intel Corp.
- 525 KARP, M. 2007. All about raid. *Netw. World News*.
- 526 MAXTOR CORPORATION. 2005a. *Atlas 15K II SAS datasheet*. Maxtor Corporation.
- 527 MAXTOR CORPORATION. 2005b. *DiamondMax 10 datasheet*. Maxtor Corporation.
- 528 NAM, Y. J., KIM, D.-W., CHOE, T.-Y., AND PARK, C. 2002. Enhancing write i/o performance of disk array  
529 rm2 tolerating double disk failures. In *Proceedings of the International Conference on Parallel Processing*  
530 *(ICPP'02)*. IEEE Computer Society, Los Alamitos, CA, 211.
- 531 PARK, C.-I. 1995. Efficient placement of parity and data to tolerate two disk failures in disk array systems.  
532 *IEEE Trans. Paral. Distrib. Syst.* 6, 11, 1177–1184.
- 533 PATTERSON, D. A., GIBSON, G., AND KATZ, R. H. 1988. A case for redundant arrays of inexpensive disks  
534 (raid). In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. ACM  
535 Press, New York, 109–116.
- 536 PLDA. 2002. *Advanced Altera Stratix PCI-X/PCI 32/64-bit Development Kit*. PLD Applications.



- 537 REED, S. AND SOLOMON, G. 1984. Polynomial codes over certain finite fields. *SIAM J. Appl. Math.* 8,  
538 300–304.
- 539 SILICON IMAGE. 2007. SiI3114 PCI to serial ATA controller datasheet. Silicon Image.
- 540 SONG, M.-A., LAN, I.-F., AND KUO, S.-Y. 2005. An area-efficient architecture of reed-solomon codec for  
541 advanced raid systems. In *Proceedings of the 11th International Conference on Parallel and Distributed*  
542 *Systems (ICPADS'05)*. IEEE Computer Society, Los Alamitos, CA, 620–626.
- 543 XILINX INC. 2006. *Hardware Accelerator for RAID6 Parity Generation/Data Recovery Controller*. Xilinx Inc.
- 544 Received October 2008; revised November 2009; accepted January 2010